

THE INVENTION CLAIMED IS:

1. A distributed data center system protocol comprising:
providing a client having a failure detector and a plurality of data centers each
including a plurality of database servers;
5 selecting one of the plurality of data centers to be a primary data center with the other
of the plurality of data centers to be a backup data center;
selecting one of the plurality of database servers in the primary data center to be a
primary database server with the other of the plurality of database servers
therein to be a backup database server and the plurality of database servers in
10 the backup data center to be backup database servers;
providing communications from the client to the primary database server and the
backup database servers;
selecting one of the backup database servers as a new primary database server when
one of the backup database servers detects a failure of the primary database
15 server; and
providing further communications from the client to the new primary database server
when the client suspects a failure of the primary database server.
2. The distributed data center system protocol as claimed in claim 1 including:
providing an abort operation from the client to the primary database server when the
20 client suspects a failure of the primary database server.
3. The distributed data center system protocol as claimed in claim 1 including:
checking whether the primary database server has already executed a transaction
operation for a transactional job corresponding to the same transactional job
before executing the transaction operation.
- 25 4. The distributed data center system protocol as claimed in claim 1 wherein:
providing the plurality of database servers includes providing local databases therefor;
and
providing the communications includes the primary database server sending a
transaction operation to the local database and executing the transaction
30 operation.
5. The distributed data center system protocol as claimed in claim 1 including:
providing a transaction operation from the client to the primary database server and
the backup database servers; and

executing the transaction operation and a second transaction operation in the same order both in the primary database server and the backup database server.

6. The distributed data center system protocol as claimed in claim 1 including: providing first and second durability levels of operation wherein employing the first durability level in the primary database server executes the transaction operation faster than the second durability level of operation.

7. The distributed data center system protocol as claimed in claim 1 wherein: providing the plurality of database servers in the primary and backup data centers includes each of the plurality of database servers having failure detectors and the client having a failure detector with properties of strong completeness and eventual weak accuracy.

8. The distributed data center system protocol as claimed in claim 1 wherein: providing the client includes providing the client with disaster detectors with properties of strong completeness and eventual strong accuracy.

9. The distributed data center system protocol as claimed in claim 1 wherein: providing the plurality of database servers in the primary and backup data centers includes each of the plurality of database servers having disaster detectors with properties of strong completeness and strong accuracy.

10. The distributed data center system protocol as claimed in claim 1 wherein: providing the primary database server includes providing a local database therefor; and

executing a transaction operation includes the primary database server sending the transaction operation to the local database.

11. The distributed data center system protocol as claimed in claim 1 including: communicating from the primary database server to the backup database servers by broadcast communication of a transaction unique identification, statements associated with the transaction operation, and control information from the primary database server to the backup database servers.

12. The distributed data center system protocol as claimed in claim 1 including: providing for only one primary data center and only one primary database server during a numbered epoch;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the backup data center detects a failure of the primary data center; and

determining the only one primary data center and the only one primary database server deterministically from the epoch number.

13. The distributed data center system protocol as claimed in claim 1 including: providing for only one primary data center and only one primary database server during a numbered epoch;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the backup data center detects a failure of the primary data center;

broadcasting the a current epoch number from the one of the plurality of database servers to the other of the plurality of database servers to change the epoch number;

determining the only one primary data center and the only one primary database server deterministically from the epoch number; and

providing a transaction operation to the only one primary data center and the only one primary database server for the epoch number.

14. The distributed data center system protocol as claimed in claim 13 including: performing a certification test in each of the plurality of database servers after receiving the broadcasting of the epoch number.

15. A distributed data center system protocol comprising: providing a client having a failure detector and a plurality of data centers each including a plurality of database servers operatively interconnected;

selecting one of the plurality of data centers to be a primary data center with the other of the plurality of data centers to be a backup data center;

selecting one of the plurality of database servers in the primary data center to be a primary database server with the other of the plurality of database servers in the primary data center to be a backup database server and the plurality of database servers in the backup data center to be backup database servers;

providing a transaction operation from the client to the primary database server and the backup database servers;

providing error messages from the backup database servers to the client indicating that the backup database servers are not the primary database server;
executing the transaction operation by the primary database server;

selecting one of the backup database servers as a new primary database server when one of the backup database servers detects a failure of the primary database server;

selecting one of the backup database servers as a new primary database server when one of the backup database servers suspects a failure of the primary database server;

providing the transaction operation from the client to the new primary database server when the client detects a failure or change of the primary database server;

executing the transaction operation by the new primary database server when the transaction operation is provided from the client to the new primary database server; and

returning the result of the executed transaction operation from the new primary database server to the client.

16. The distributed data center system protocol as claimed in claim 15 including: providing an abort operation from the client to the primary database server when the client detects a failure of the primary database server.

17. The distributed data center system protocol as claimed in claim 15 including: checking whether the primary database server has already executed a transaction operation for a transactional job corresponding to the same transactional job before executing the transaction operation.

18. The distributed data center system protocol as claimed in claim 15 wherein: providing the plurality of database servers includes providing local databases therefor; executing the transaction operation includes the primary database server sending the transaction operation to the local database and waiting for a reply;
executing the transaction operation includes waiting for a reply from the local database to the primary database server and sending the reply to the client;
executing the transaction operation includes sending a commit request to the primary database server in response to the reply; and

executing the transaction operation includes the primary database server broadcasting a transaction unique identification, SQL statements associated with the transaction operation, and control information to the backup data server in the primary data center, and the plurality of backup data servers in the backup data center.

19. The distributed data center system protocol as claimed in claim 15 including: providing a second transaction operation from the client to the primary database server and the backup database servers; and

executing the transaction operation and the second transaction operation in the same order both in the primary database server and the backup database server.

20. The distributed data center system protocol as claimed in claim 15 including: providing first and second durability levels of operation wherein:

employing the first durability level in the primary database server executes the transaction operation faster than the second durability level of operation and

employing the second durability level in the primary database server executes the transaction operation with the assurance that if the primary data center suffers a disaster, the plurality of backup databases in the backup data center will receive the transaction operation.

21. The distributed data center system protocol as claimed in claim 15 wherein:

providing the plurality of database servers in the primary and backup data centers includes each of the plurality of database servers having failure detectors and the client having a failure detector with the properties of:

strong completeness wherein, if the primary database server fails at time t , then there is a time $t' > t$ after which the primary database server is permanently suspected of failure by the client and by the backup database server; and

eventual weak accuracy wherein, if the primary database server that does not fail, then there is a time after which the primary database server is never suspected of failure by the client and by the backup database server.

22. The distributed data center system protocol as claimed in claim 15 wherein:

providing the client includes providing the client with disaster detectors with the properties of:

strong completeness wherein, if the primary data center fails at time t , then there is a time $t' > t$ after which the primary data center is permanently suspected of failure by the client; and

eventual strong accuracy wherein, if the primary data center that does not fail, then there is a time after which the primary data center is never suspected of failure by the client.

23. The distributed data center system protocol as claimed in claim 15 wherein: providing the plurality of database servers in the primary and backup data centers includes each of the plurality of database servers having disaster detectors with the properties of:

strong completeness wherein, if the primary data center fails at time t , then there is a time $t' > t$ after which the primary data center is permanently suspected of failure by the backup database servers; and

strong accuracy wherein, if the primary data center that does not fail, then the primary data center is never suspected of failure by the backup database servers.

24. The distributed data center system protocol as claimed in claim 15 wherein: providing the primary database server includes providing a local database therefor; and

executing the transaction operation includes the primary database server sending the transaction operation to the local database.

25. The distributed data center system protocol as claimed in claim 15 including: communicating from the primary database server to the backup database servers by broadcast communication of a transaction unique identification, statements associated with the transaction operation, and control information from the primary database server to the backup database servers.

26. The distributed data center system protocol as claimed in claim 15 including: providing for only one primary data center and only one primary database server during a numbered epoch;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the backup data center detects a failure of the primary data center;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the primary or backup data centers detects a failure of the primary database server;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the primary or backup data center suspects a failure of the primary database server;

determining the only one primary data center and the only one primary database server deterministically from the epoch number; and

providing the transaction operation to the only one primary data center and the only one primary database server for the epoch number.

27. The distributed data center system protocol as claimed in claim 15 including: providing for only one primary data center and only one primary database server during a numbered epoch;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the backup data center detects a failure of the primary data center;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the primary or backup data centers detects a failure of the primary database server;

establishing an epoch number by one of the plurality of database servers when the one of the plurality of database servers in the primary or backup data center suspects a failure of the primary database server;

broadcasting the a current epoch number from the one of the plurality of database servers to the other of the plurality of database servers to change the epoch number;

determining the only one primary data center and the only one primary database server deterministically from the epoch number; and

providing the transaction operation to the only one primary data center and the only one primary database server for the epoch number.

28. The distributed data center system protocol as claimed in claim 27 including: performing a certification test in each of the plurality of database servers after receiving the broadcasting of the epoch number wherein:

executing the transaction operation includes broadcasting the transaction operation and the epoch number from the primary database server; and committing the completed transaction operation when the certification test determines the epoch number in the new primary database server is the same as the epoch number of the transaction operation.

5